

Unsupervised Disentanglement of Linear-Encoded Facial Semantics

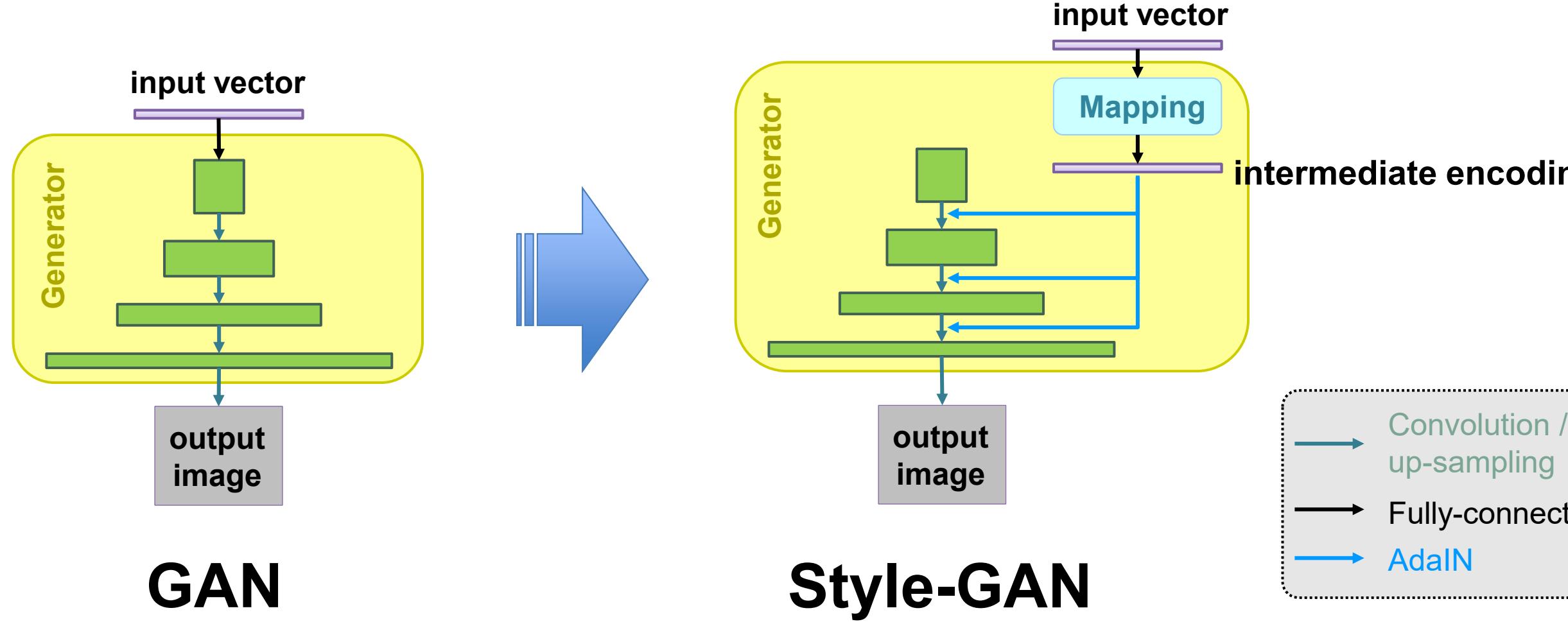
Yutong Zheng, Yu-Kai Huang, Ran Tao, Zhiqiang Shen, and Marios Savvides

{yutongzh,yukaih2,rant,zhiqians,mariooss}@andrew.cmu.edu



Background

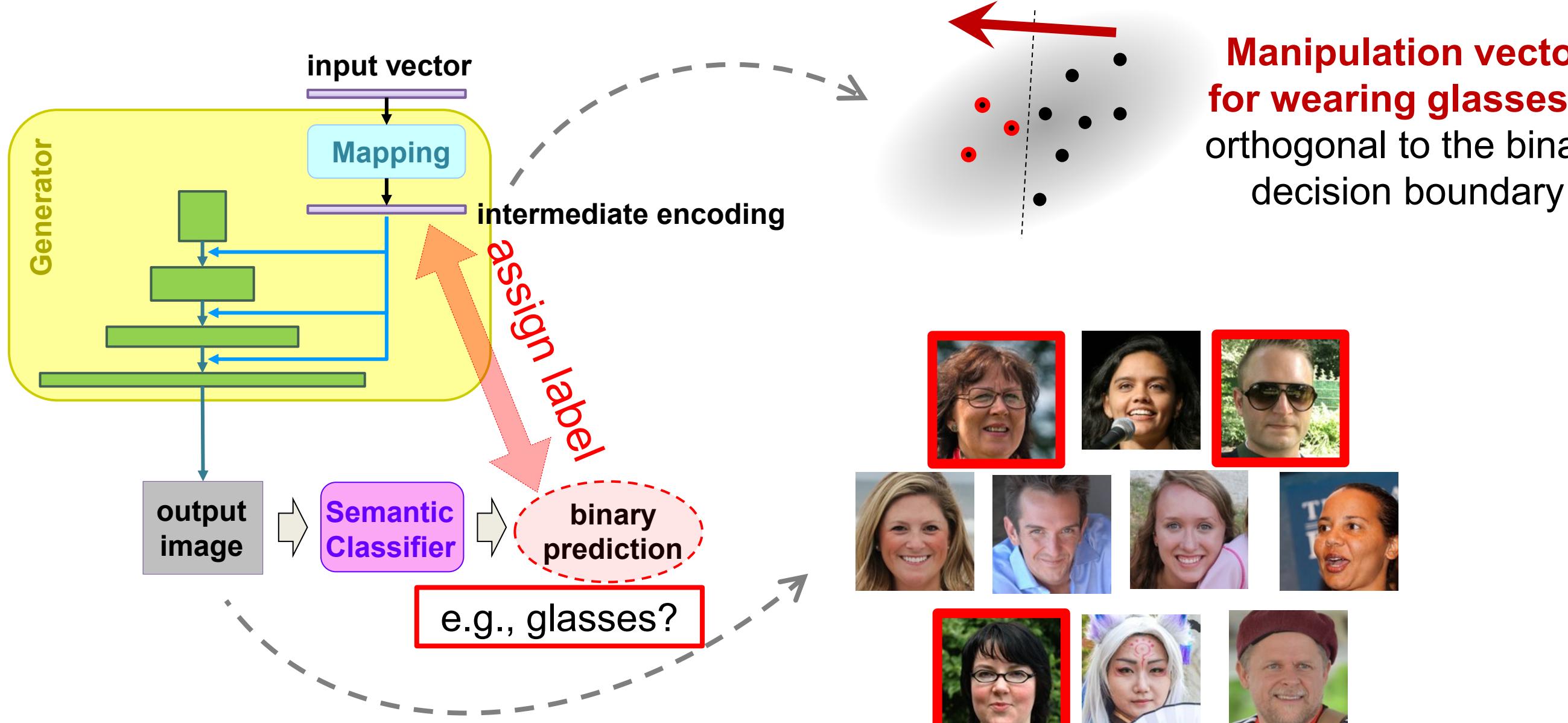
Style-GAN makes representations linear-encoded:



Intermediate encoding with AdalN operation provides more flexibility to the representations, making linear encoding possible.

Previous Work: Manipulation Vectors

Finding manipulation vectors with pretrained semantic classifiers:

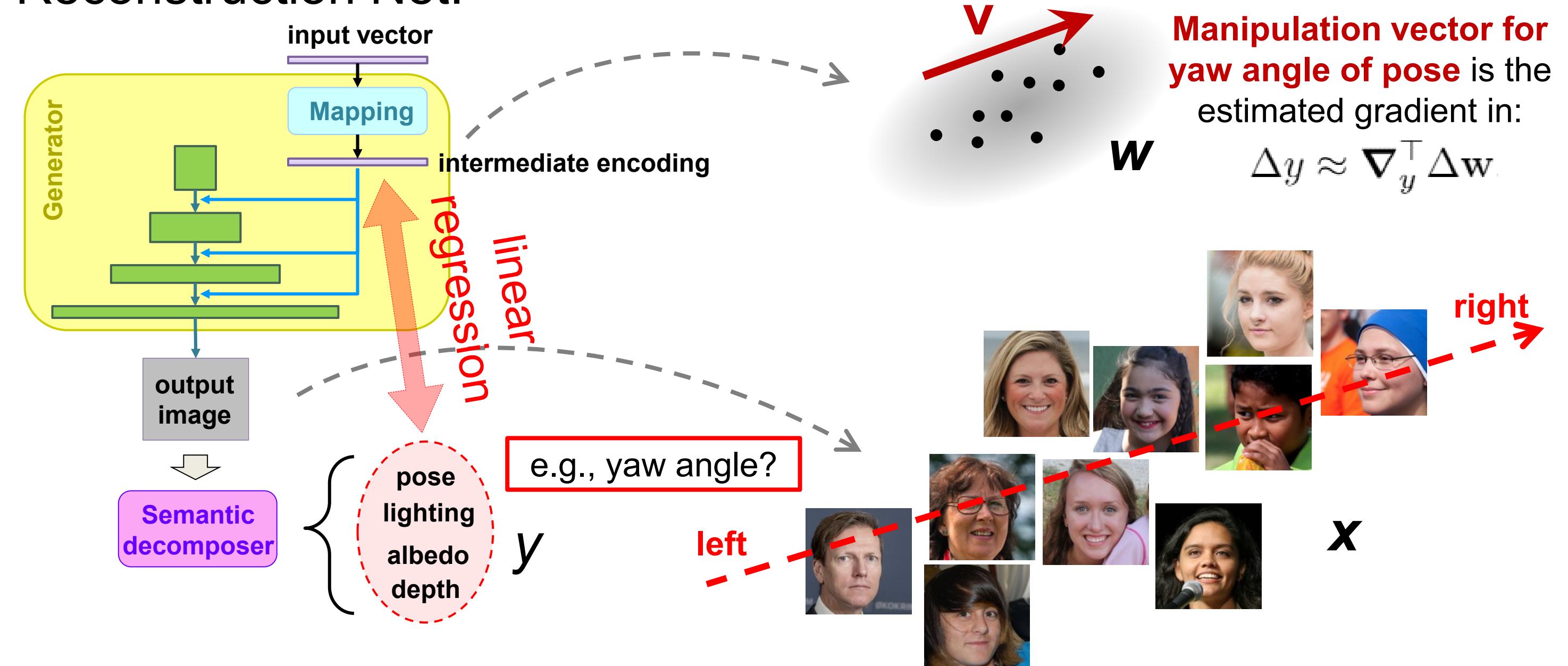


Drawbacks for *supervised* methods:

- Requires pre-training **semantic classifiers** with labeled data.
- Semantics are artificially defined, resulting in potential redundancy and incompleteness.

The Method: Going Unsupervised

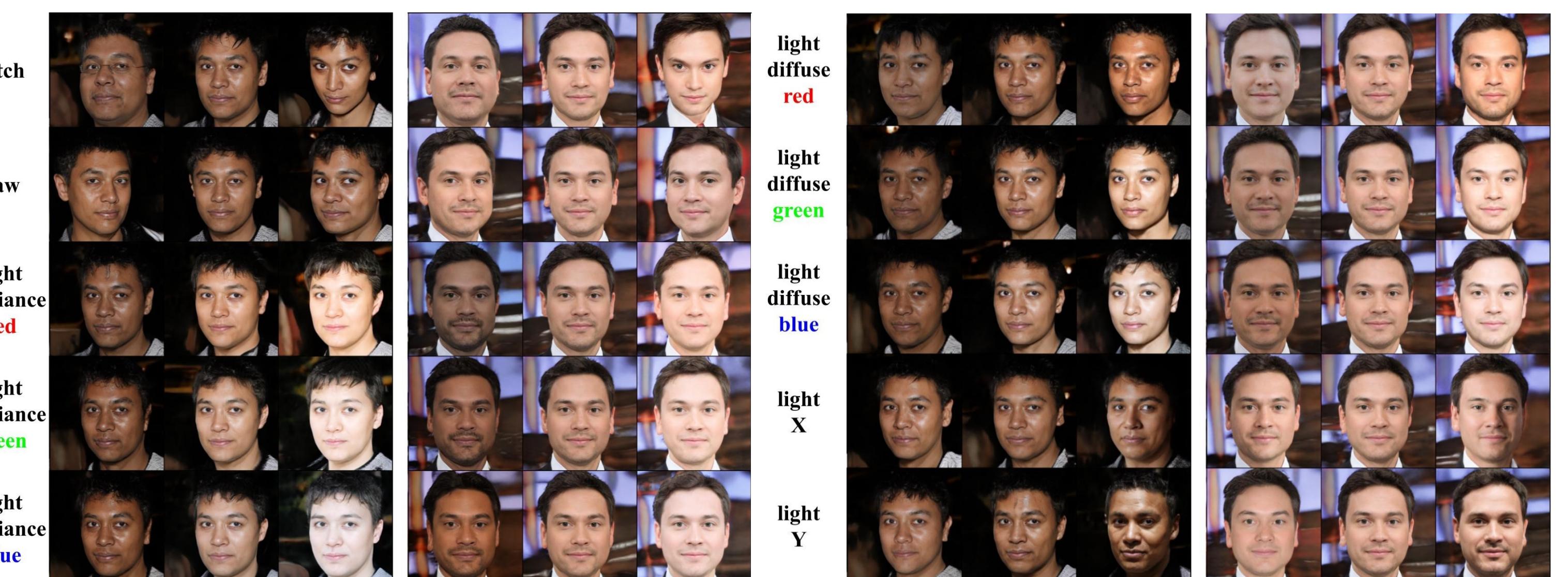
Setup: Unsupervised Pre-training of **Style-GAN generator** and **semantic decomposer** from the 3D Deformable Face Reconstruction Net.



Gradient (v) for any scalar semantic components (y): sample a bunch of w and output y , then optimize the objective:

$$\min_v \|\Delta Y - \Delta W v\|_2^2$$

Manipulate v associated with interpretable scalar semantic components: pose and lighting:

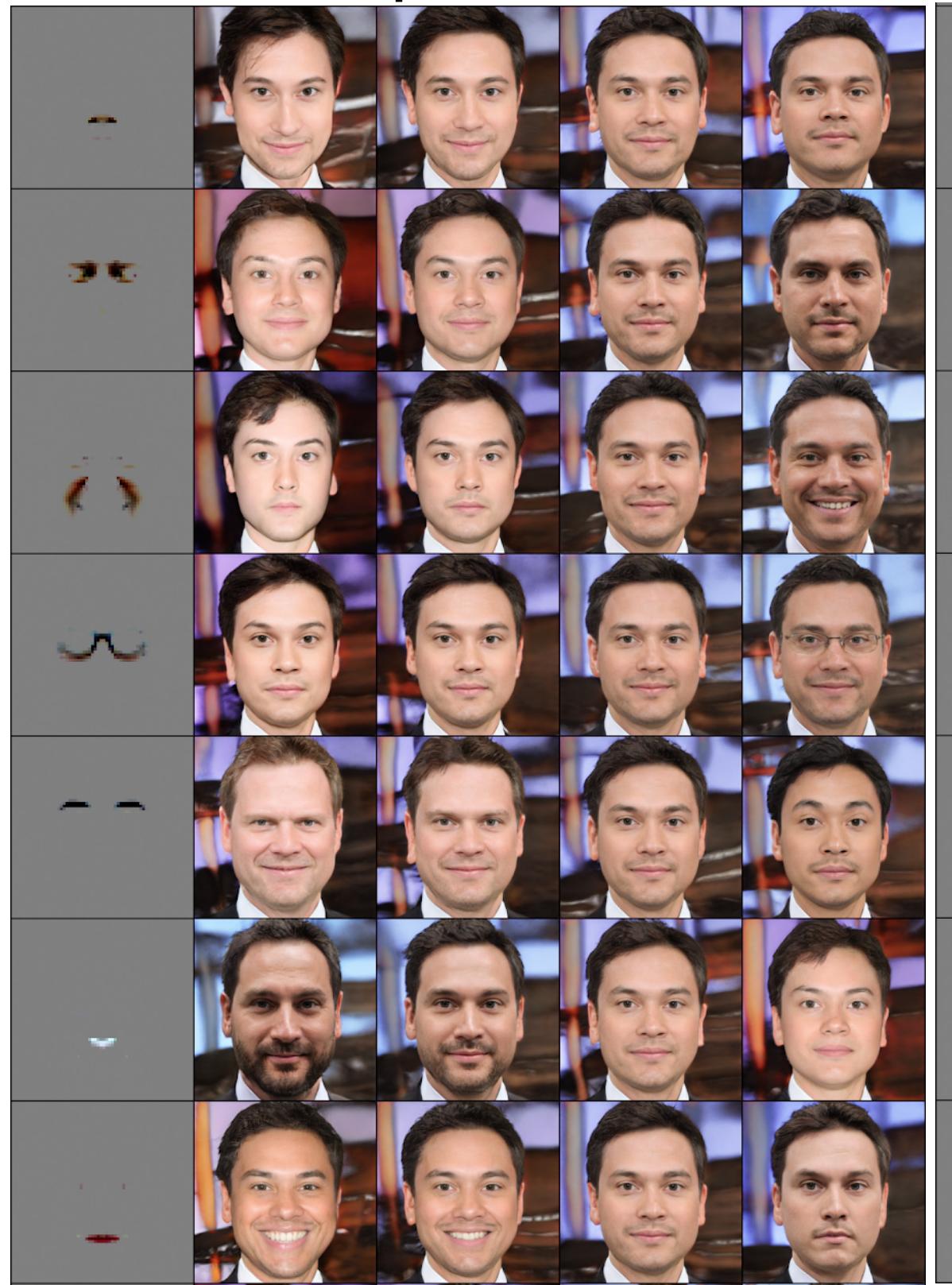


The Method: Localized Representations

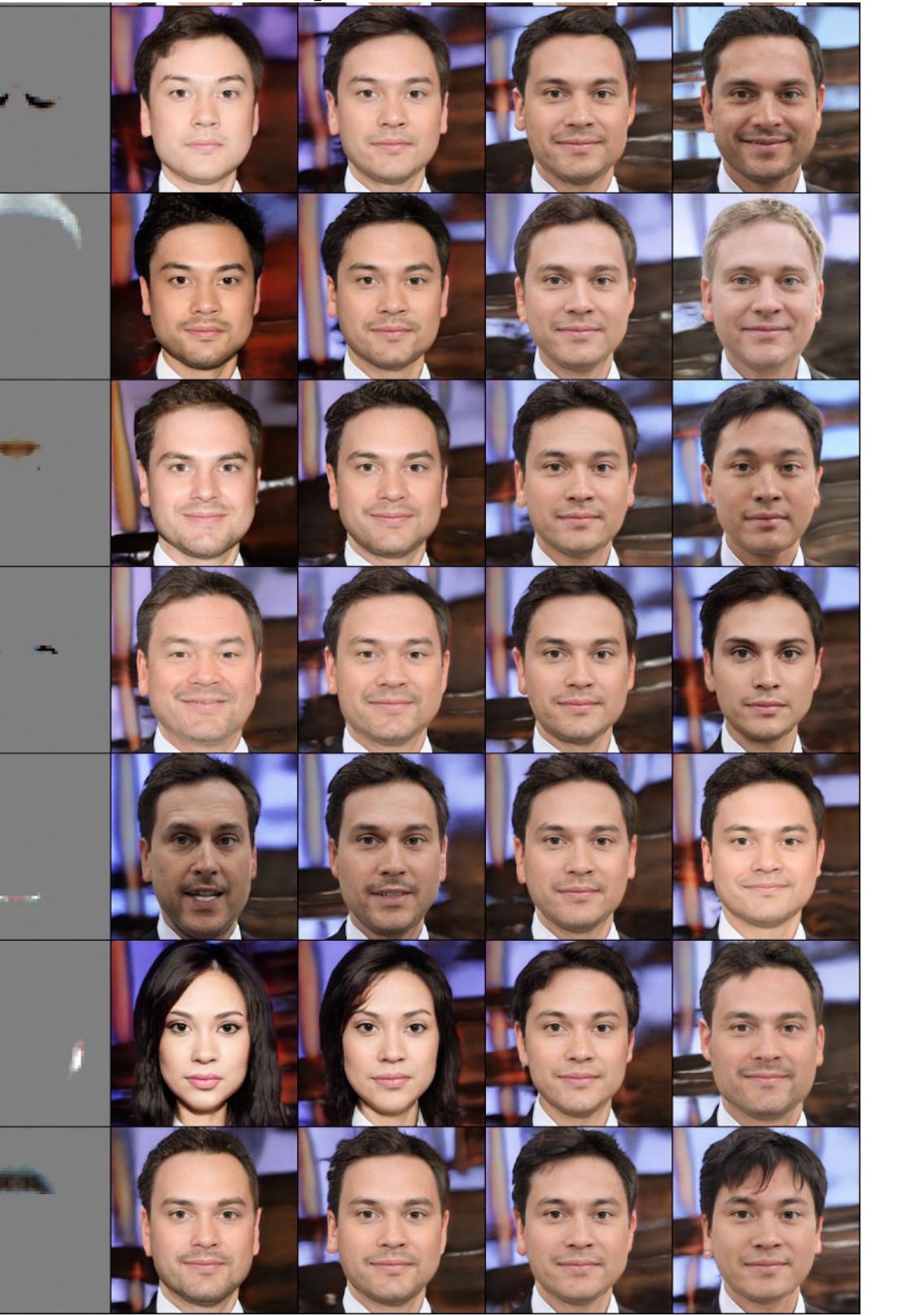
Gradient (\tilde{v}) for pixel-level semantic components: combine v in albedo and depth from previous step as a Jacobian J_v^* , then estimate local facial variations by optimizing:

$$\begin{aligned} \min_{U, \tilde{V}} & \|J_v^* - U \tilde{V}^\top\|_F^2 + \alpha \|U\|_1 + \beta \sum_{i \neq j} (\hat{v}_i^\top \hat{v}_j)^2 \\ \text{s.t. } & \|\hat{v}_p\|_2 = 1, \end{aligned}$$

U Manipulation results



U Manipulation results



Different components in \tilde{v}

Cosine distances of \tilde{v} comply with the contextual relations of local semantics:

